

Lab 7: How the Tetracycline Came to Peoria

36-350

10 October 2014

Agenda: Transforming data; combining information from multiple objects; practice with selective access; practice applying functions.

Now-common ideas like “early adopters” and “viral marketing” grew from sociological studies of the diffusion of innovations. One of the most famous of these studies tracked how a then-new antibiotic, tetracycline, spread among doctors in four small cities in Illinois in the 1950s. In this lab, we will go back to that data to look at one of the crucial ideas, that of the innovation (prescribing tetracycline) “spreading” from person to person.

On the class website, you will find two data files, [http://www.stat.cmu.edu/~cshalizi/statcomp/14/labs/07/ckm_nodes.csv] and [http://www.stat.cmu.edu/~cshalizi/statcomp/14/labs/07/ckm_network.dat]. The former has information about each individual doctor in the four towns. The latter records which doctors knew each other.

Part I

1. Load the dataset `ckm_nodes.csv` into a data frame, `ckm_nodes`. Check that it has 246 rows and 13 columns. Check that there are columns named `city` and `adoption_date`.
2. `adoption_date` records the month in which the doctor began prescribing tetracycline, counting from November 1953. If the doctor did not begin prescribing it by month 17, i.e., February 1955, when the study ended, this is recorded as `Inf`. If it’s not known when or if a doctor adopted tetracycline, their value is `NA`.
 - a. How many doctors *began* prescribing tetracycline in each month of the study? How many never prescribed? How many are NAs? *Hints:* `table()`, `is.na()`, `sum()`.
 - b. Create a vector which records the index numbers of doctors for whom `adoption_date` is not `NA`. Check that this vector has length 125. Re-assign `ckm_nodes` so it only contains those rows. (Do not drop rows if they have a value for `adoption_date` but are `NA` in some other column.) Use this cleaned version of `ckm_nodes` for the rest of the lab.
3. Create plots of the number of doctors who began prescribing tetracycline each month versus time. (It is OK for the numbers on the horizontal axis to just be integers rather than formatted dates.) Produce another plot of the *total* number of doctors prescribing tetracycline in each month. The curve for total adoptions should first rise rapidly and then level out around month 6.
4. *Adopted already or not yet?*
 - a. Create a Boolean vector which indicates, for each doctor, whether they had begun prescribing tetracycline by month 2. Convert it to a vector of index numbers. There should be twenty such doctors.
 - b. Create a Boolean vector which indicates, for each doctor, whether they began prescribing tetracycline after month 14, or never prescribed it. Convert it to a vector of index numbers. There should be twenty-three such doctors.

Part II

5. The file `ckm_network.dat` contains a binary matrix; the entry in row i , column j is 1 if doctor number i said that doctor j was a friend or close professional contact, and 0 otherwise. Load the file into R as `ckm_network`, and verify that gives you a square matrix which contains only 0s and 1s, and that it has 246 rows and columns. Drop the rows and columns corresponding to doctors with missing `adoption_date` values. Check that the result has 125 rows and columns. Use this reduced matrix, and its row and column numbers, for the rest of the lab.
6. Create a vector which stores the number of contacts each doctor has. Do not use a loop. Check that doctor number 41 had 3 contacts.
Hint: You could do this using `apply`, but you can also do it in one line with a single function.
7. *Counting Peer Pressure*
 - a. Create a Boolean vector which indicates, for each doctor, whether they were contacts of doctor number 37, *and* had begun prescribing tetracycline by month 5. Count the number of such doctors without converting the Boolean vector to a vector of indices. There should be three such doctors.
 - b. What proportion of doctor 37's friends do those two doctors represent?

We will continue with this data set in the next lab.

Behind the scenes: The original study was published as

James Coleman, Elihu Katz and Herbert Menzel, "The Diffusion of an Innovation Among Physicians" *Sociometry* **20** (1957): 253–270.

The files used here are taken from [<http://moreno.ss.uci.edu/data.html#ckm>] with some formatting changes. CKM actually measured three types of link among the doctors — friendship, general discussion, and asking for medical advice. To keep things simple, we are combining all three types of tie, and treating them as symmetric.